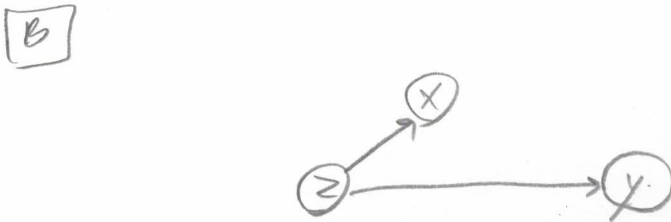


MATH 303 LEC 19
 Now let's put x into the DAG. For the setting of prediction (292) and usually in inference (391) eg. in experimental studies, the x 's must be measured before the response is measured. Let's look at some ^{very} simple cases



Here x does not cause y and x is not correlated to y .

However, a spurious correlation may occur. A model $g(x)$ would perform worse than g_0 since it is overfit (est. err).



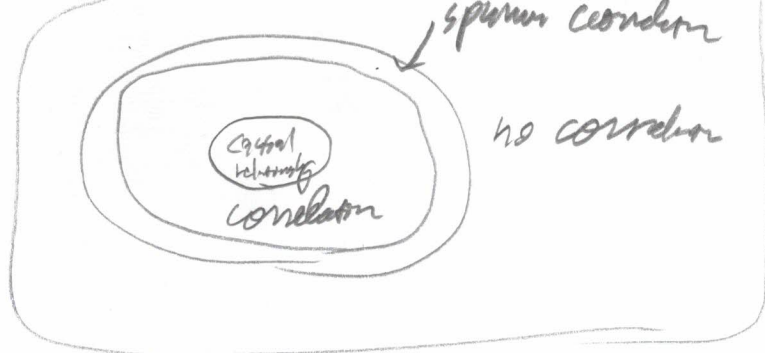
Here x does not cause y but x is correlated to y .

Here, $g(x)$ should do better than g_0 depending on how strong the relationship $x = k_{x,z}(z) + \text{noise}$ is and how strong $y = k_{y,z}(z) + \text{noise}$ is. Carefully it is captured when the model $y \sim x$ is constructed



x both causes y and x is corr. to y . $g(x)$ should do better than g_0 depending on how strong the relationship is in $z = k_{z,x}(x) + \text{noise}$ and how carefully it is captured when model $y \sim x$ is constructed and $y = k_{y,x}(x) + \text{noise}$

Thus the full picture,
for $\forall x, y \dots$
the relationship between



Thus correlation $\not\Rightarrow$ causation, but sometimes they coincide.
How to prove causation? That's the big topic to be explored.

Thus causation \Rightarrow correlation.

Correlation \Rightarrow causation sometimes is through
a common effect (which is called a
confounder) or "lurking variable".

In [5], z is the "confounder" of the relationship of x on y .
If I have more data, the definition will be more exact.

What does the real picture look like with multiple x 's?

$$y = f(z_1, \dots, z_q)$$

$$y = f(x_1, \dots, x_p) + \epsilon$$

Each scenario

A, B, C are

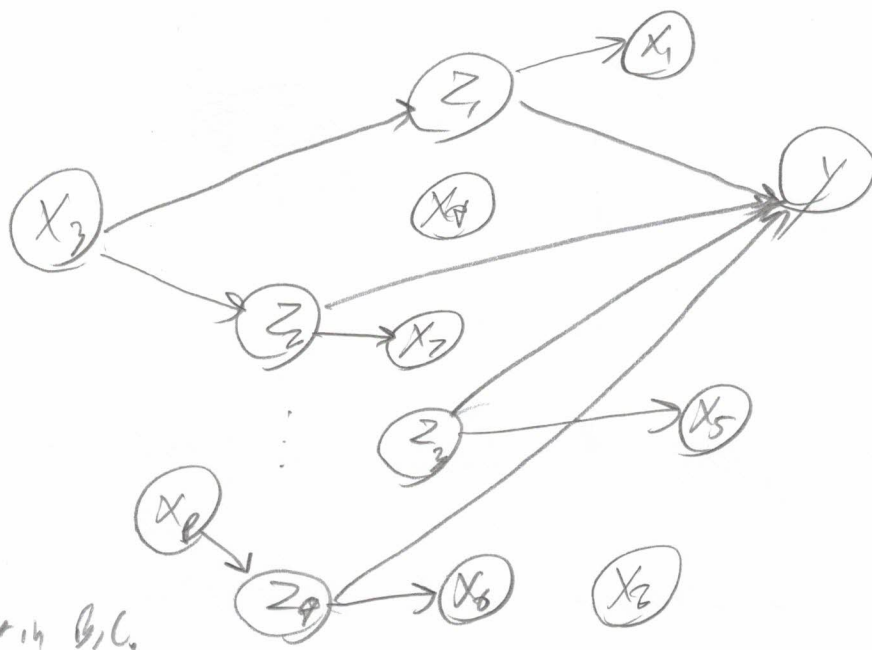
represented and more!

You don't know precisely

in A (they're useless

and will cause

overfitting). But you know in B, C .



the answer of

Given this very messy phrase, what does $\vec{\beta}$ mean in
 OLS, logistic equation, poisson regression, weibull, cox ph, etc?
 let's be careful and build up an interpretation iteratively

OLS /
 robust
 regression

$$\hat{y} = b_0 + b_1 x_1 + \dots + b_p x_p$$

bernoulli
 regression

e.g. $\phi^{-1}(\hat{p}) = \ln\left(\frac{\hat{p}}{1-\hat{p}}\right)$ for logistic link

$$\phi^{-1}(\hat{p}) = b_0 + b_1 x_1 + \dots + b_p x_p$$

Poisson
 regression

$$\hat{y} = e^{b_0} e^{b_1 x_1} \dots e^{b_p x_p} \Rightarrow \ln(\hat{y}) = b_0 + b_1 x_1 + \dots + b_p x_p$$

Neg binom
 regression

$$\hat{y} = \Gamma\left(1 + \frac{1}{\hat{\mu}}\right) e^{b_0} e^{b_1 x_1} \dots e^{b_p x_p} \Rightarrow \ln(\hat{y}) = \ln\left(1 + \frac{1}{\hat{\mu}}\right) + b_0 + b_1 x_1 + \dots + b_p x_p$$

Weibull
 regression

$$h(t) = h_0(t) e^{b_0} e^{b_1 x_1} \dots e^{b_p x_p} \Rightarrow \ln(h(t)) = \ln(h_0(t)) + b_0 + b_1 x_1 + \dots + b_p x_p$$

Cox PH
 regression

What, consider $j=1$, the first covariate, and let it be "blood sugar". let
 it be measured in mg/dL.

Here's the causal interpretation for OLS for x_j = blood sugar measured in mg/dL and y = blood triglycerides measured in mmol/L

"If blood sugar is increased by one unit of metric x_j and all other
name of metric x_j

for $b_j = 1.17$
 and $SE_j = 0.25$

parameters remain constant, blood triglycerides will realistically increase
name of linear response "name" of $b_j > 0$ or "decrease"

by an estimated 1.17 ± 0.25 mmol/L
 value of $b_j \pm SE_j$ unit of response y

assuming the blood triglycerides linear in the P
the response function assumption

Focus is on the manipulation / interpretation of x_j

Let's do this for logistic regression where $x_j = \text{blood sugar (mg/dL)}$ and $y = 1$ if person will get diabetes. and $b_j = 0.49$, $s_{b_j} = 0.12$

"If blood sugar is increased by one mg/dL and all other measurements remain constant, the log odds of getting diabetes will roughly increase by an amount 0.49 ± 0.12 assuming the log odds of getting diabetes is linear in these p covariates."

How about for weibull? $x_j = \text{blood sugar (mg/dL)}$ and $y = \text{risk survival (yr)}$ and $b_j = -0.02$, $s_{b_j} = 0.007$

"If blood sugar is increased by one mg/dL and all other measurements remain constant, the log survival will roughly decrease by an amount 0.02 ± 0.007 assuming the survival is weibull distributed with log mean linear in these p covariates."

Alternatively, if he wants to use survival in years, he can use the proportional hazards regression. $e^{b_j} = e^{-0.02} = 0.98 \Rightarrow -2.0\%$. For s_{b_j} , more difficult, real delta method

"If blood sugar is increased by one mg/dL and all other measurements remain constant, the survival will roughly decrease by an amount 2.0% assuming"
 prop. hazard regression

For cox PH, same covariates

"If blood sugar is increased by one mg/dL and all other measurements remain constant, the log hazard rate of survival will roughly decrease by an amount 0.02 ± 0.007 assuming the survival process has a hazard rate log linear in the p covariates."